

All You Ever Need to Know About Counterfactual Explanations

(and were afraid to ask...)

Fundamentals, Methods, & User Studies for XAI

André Artelt ~ Ulrike Kuhl ~ Mark T. Keane





André Artelt

Bielefeld University, Germany
University of Cyprus, Cyprus



UNIVERSITÄT
BIELEFELD



Faculty of Technology



University
of Cyprus



Ulrike Kuhl

Faculty of Technology
Bielefeld University
Bielefeld, Germany



UNIVERSITÄT
BIELEFELD



Faculty of Technology

dataninja.nrw

Funded by the
Ministry of Culture and Science
of the State of
North Rhine-Westphalia



Mark T. Keane

School of Computer Science
University College Dublin
Dublin, Ireland



University College Dublin
Ireland's Global University

Insight

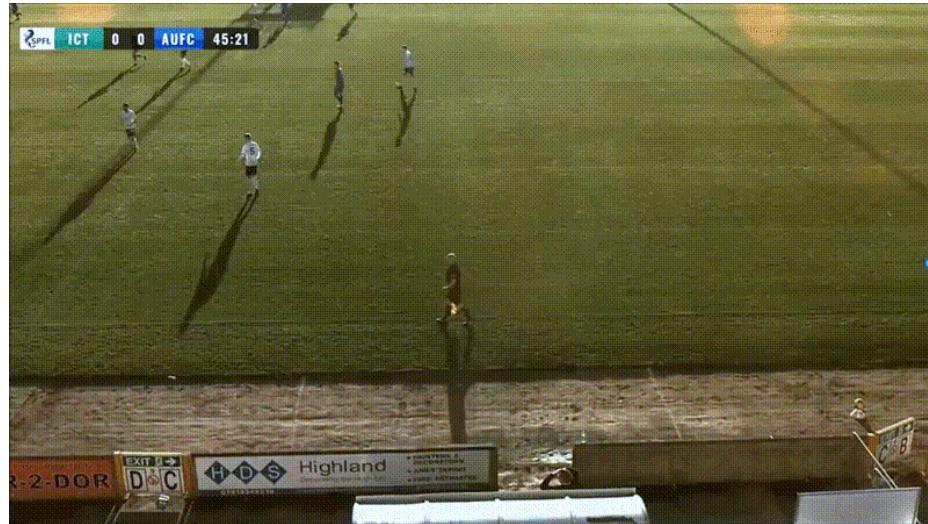
SFI RESEARCH CENTRE FOR DATA ANALYTICS

The Problem

ChatGPT OpenAI



<https://youtu.be/eMx-2s7mZ24>

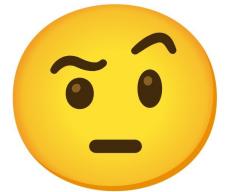


AI is out there...being deployed...busily making mistakes !

Spot the Ball !



Spot the Head !



One Solution: Explain It !

- *Ante-Hoc*: Somehow show the workings of the machine (*simulability*)
- *Post-Hoc*: Justification after-the-fact to *explain* workings (*evidence of what happened*)



Post Hoc: It's Wrong Because !

Test Image



Label: **Rifle**
Classification: **Shovel**

Explanation-by-Example



Label: **Shovel**
Classification: **Shovel**

Post Hoc: To Audit a CNN !



SPECIAL SECTION ON DATA-ENABLED INTELLIGENCE FOR DIGITAL HEALTH

Received September 14, 2019, accepted October 7, 2019, date of publication October 21, 2019, date of current version October 31, 2019.

Digital Object Identifier 10.1109/ACCESS.2019.2948430

Explainable Prediction of Chronic Renal Disease in the Colombian Population Using Neural Networks and Case-Based Reasoning

GABRIEL R. VÁSQUEZ-MORALES¹, SERGIO M. MARTÍNEZ-MONTERUBIO²,

PABLO MORENO-GER³, AND JUAN A. RECIO-GARCÍA¹

¹Office of Information and Communications Technology, Ministry of Health and Social Protection, Bogotá 110311, Colombia

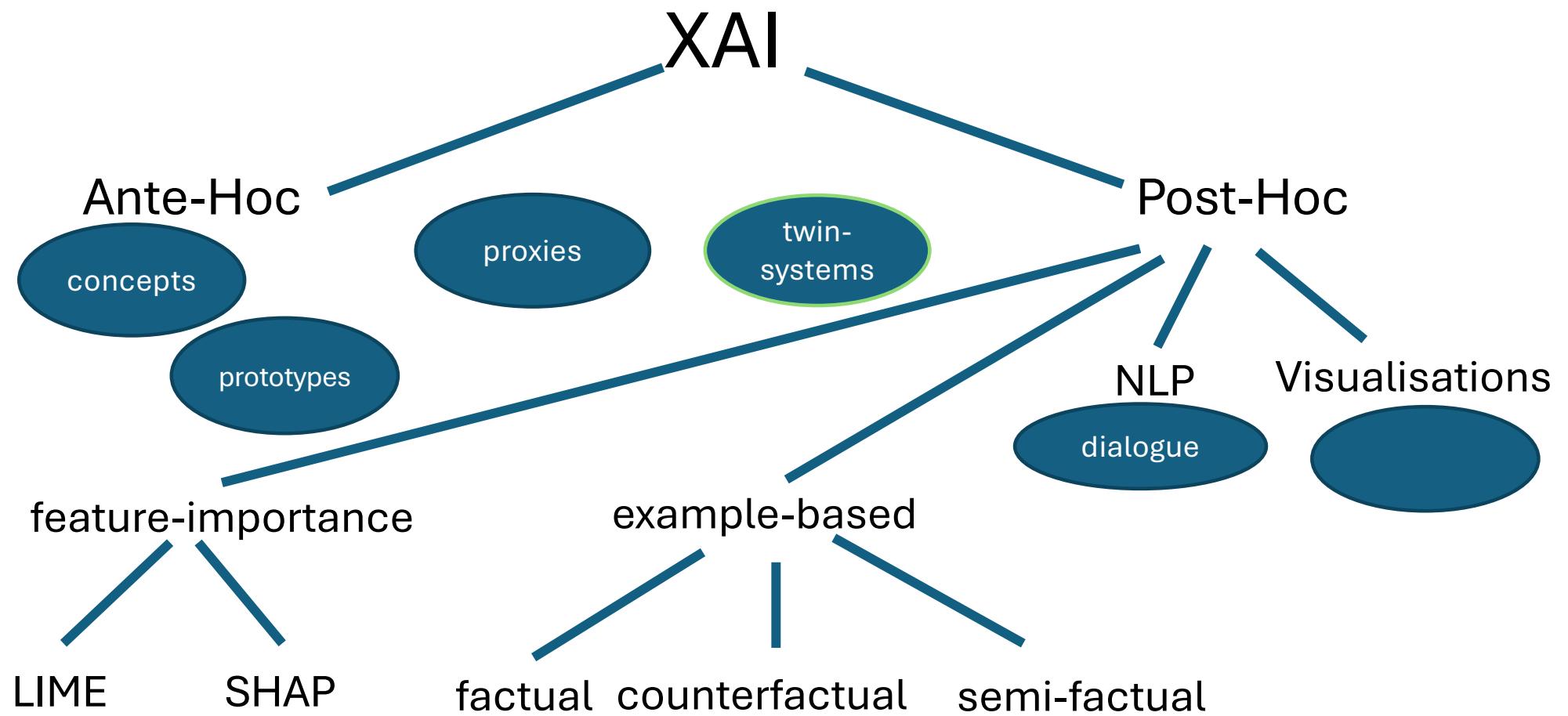
²Group of Artificial Intelligence Applications, Department of Software Engineering and Artificial Intelligence, Faculty of Computer Science, Universidad Complutense de Madrid, Ciudad Universitaria, 28040 Madrid, Spain

³School of Engineering, Universidad Internacional de La Rioja (UNIR), 26006 Logroño, Spain

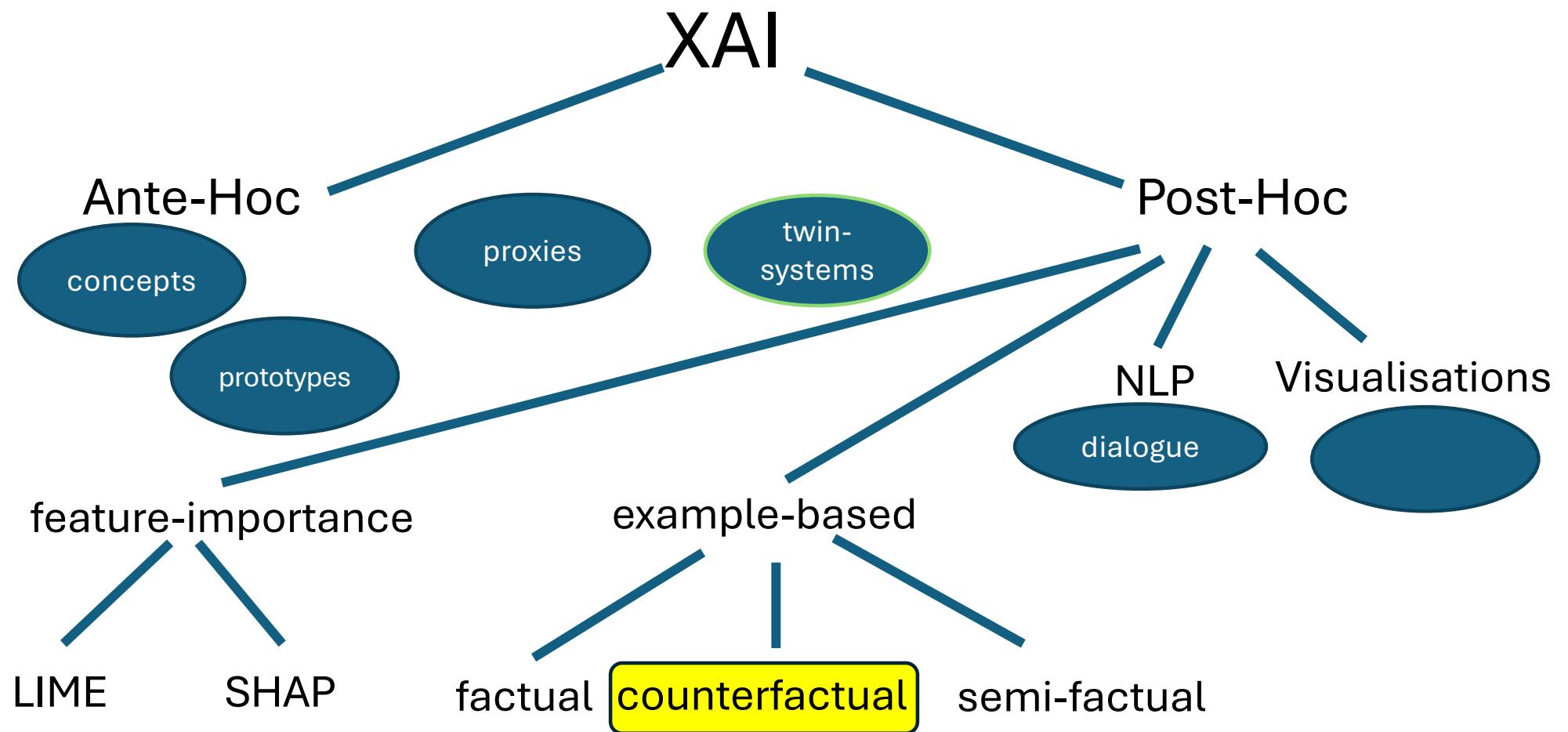


FIGURE 17. Map of coverage of the population at risk of developing CKD.

Explainable AI: Taxonomy



Explainable AI: Taxonomy

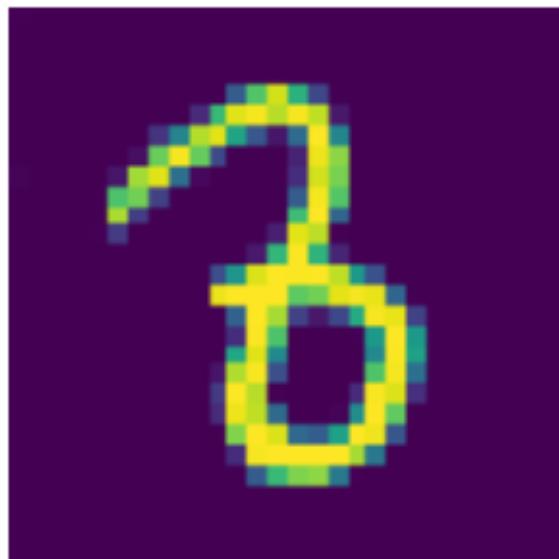


Counterfactual Explanations



Counterfactual Explanations

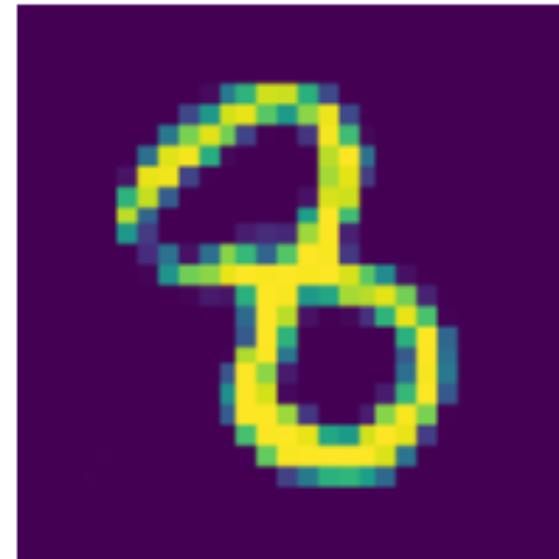
If the test image looked like this, I would have thought it was an “8”.



Test Image

Label: 8

Prediction: 3

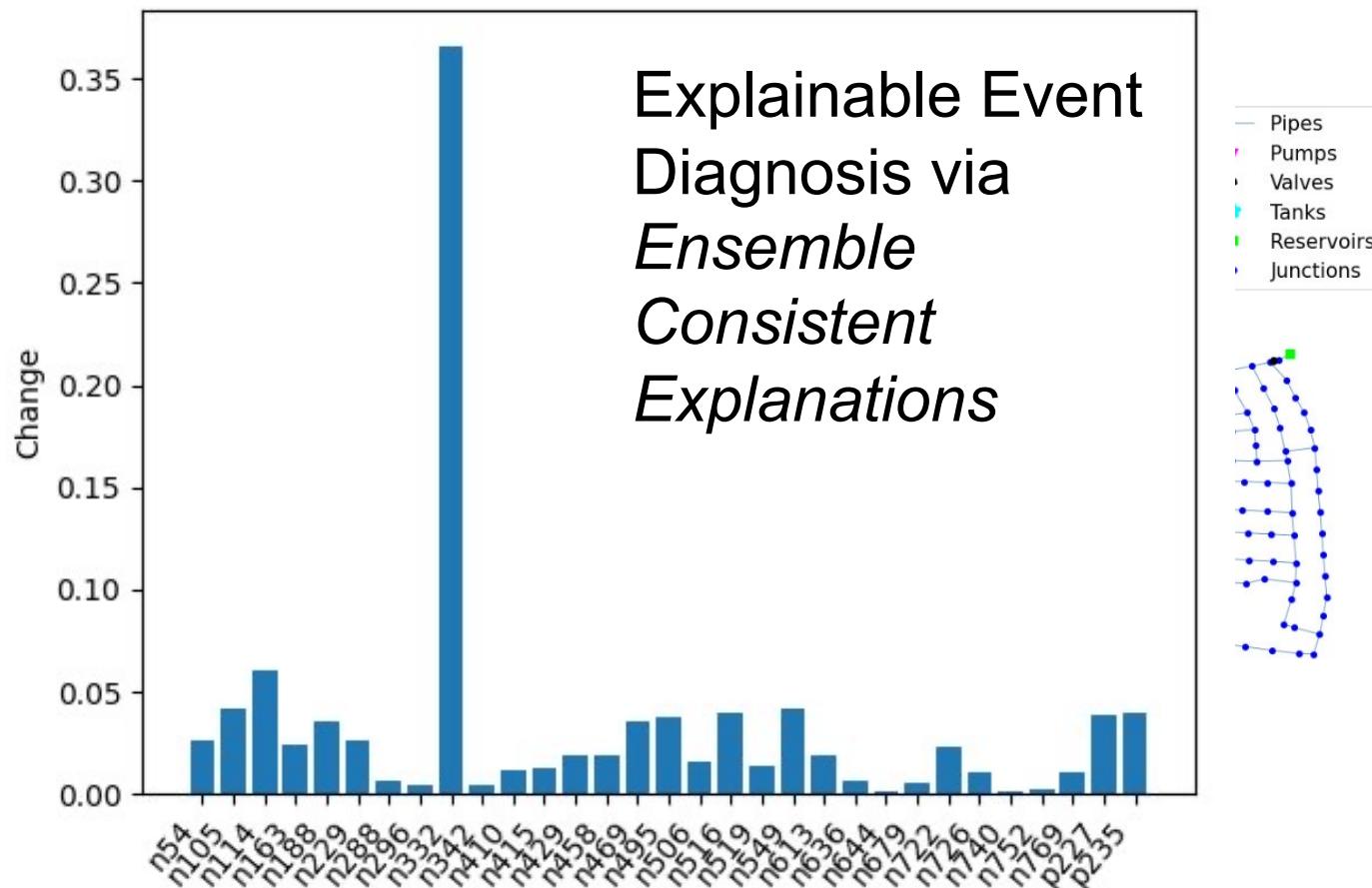


Counterfactual

New Prediction: 8

Counterfactual Explanations

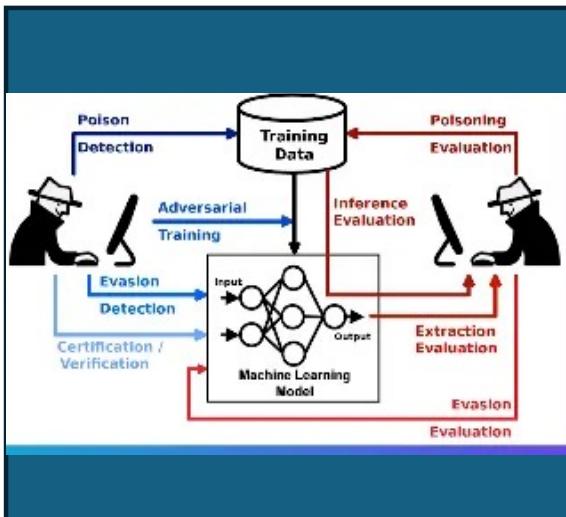
Explaining Anomalies Across Virtual Sensors



[Artelt et al. 2022 XAI Workshop @ IJCAI](#)

Why Counterfactuals?

Artificial Intelligence



Looks like adversarial learning...

[Wachter et al \(2019\) Faccit-19](#)

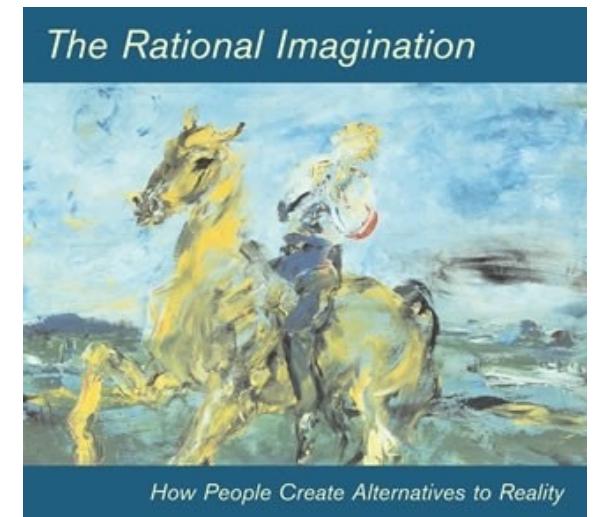
Legal



GDPR compliant...

[Wachter et al \(2017\) Harv. JL &B Tech.](#)

Psychology



Used by people spontaneously...

[Byrne \(2007\) MIT Press](#)

Plan for the Day...

CF Tutorial	
TIME	Topics
9:00 AM	Introduction <i>Hello and Introducing Ourselves!</i> <i>Hands-on: Trying Our Study (follow link)</i>
9:30 AM	Historical Fundamentals of Counterfactuals <i>From Philosophy to XAI (via Psychology)</i> <i>Two Sample User Studies and Q&A</i>
10:30 AM	COFFEE (10:30-11:00)
11:00 AM	Fundamentals of Counterfactuals in AI <i>Formalisation</i> <i>Modelling Approaches & Key Constraints</i>
11:30 AM	Using Counterfactual Algorithms <i>Hands-on: A Counterfactual Toolbox (AA)</i> <i>Hands-on: Checking Out Notebooks and Q&A</i>
12:00 PM	Fundamentals of User Studies <i>User Studies I: A Simple Two-Group Design</i>
12:30 PM	LUNCH (12:30-14:00)
2:00 PM	Algorithmic Growth Points <i>Computational Future Directions and Q&A</i>
2:30 PM	More Fundamentals of User Studies <i>User Studies I: A Simple Two-Group Design (cont.)</i>
3:00 PM	COFFEE (15:00-15:30)
3:30 PM	From Fundamentals to an Actual User Study <i>User Studies II: A More Complex Design</i> <i>User Studies III: Even More Complex Designs</i> <i>Hands-on: Looking At Our Study</i>
5:00 PM	Closing Session, Discussion and Final Q&A TUTORIAL END



**The Full Experience:
Be a Participant!**



<https://tinyurl.com/ijcai24-XAI>

